

SPM project: Parallel Prefix

Gaspere Ferraro, 520549

Master Degree in Computer Science - University of Pisa

ferraro@gaspa.re

November 3, 2018

1 Introduction

In this report we will analyze, theoretically and practically, the resolution of the problem of the (parallel) prefix sum:

Given a vector $x = \langle x_0, x_1, \dots, x_{n-1} \rangle$ and a binary operation \oplus compute the vector $y = \langle x_0, x_0 \oplus x_1, x_0 \oplus x_1 \oplus x_2, \dots, x_0 \oplus x_1 \oplus \dots \oplus x_{n-1} \rangle$.

In literature this operation is also called (inclusive) scan or partial sum.

For the analysis of the problem we have to make two assumptions:

- The binary operation \oplus is associative ($a \oplus (b \oplus c) = (a \oplus b) \oplus c$) and commutative ($a \oplus b = b \oplus a$), this is an important assumption as we will see in the next chapters the order of the operations may not be preserved.
- The size of the input vector is a power of 2, not a strong assumption, as all the algorithms we will present could be easily generalize to all the sizes, but it only helps to simplify some operations.

2 Sequential algorithm

The sequential algorithm simply compute each element of the vector y as follow:

$$y_i = \begin{cases} x_0 & \text{for } i = 0 \\ x_i \oplus y_{i-1}, & \text{for } 0 < i < n \end{cases}$$

This algorithm is optimal in a sequential model as it has a running time of $\mathcal{O}(n)$, assuming that \oplus is $\mathcal{O}(1)$, and performs $n - 1$ calls to \oplus operation.

3 Parallel architecture design

The problem of computing the prefix sum vector is a classical example of a problem that have an optimal solution in a sequential model but that can be optimized in a parallel model. The optimization is not in terms of total complexity or in the number of \oplus operations performed (which are already optimal in the sequential algorithm) but in terms of completion time. When the number of available workers is more than one we can trade-off more total work for less completion time.

We will now introduce two different algorithms that solve in an efficient way the prefix sum problem in a parallel model.

3.1 Block-based algorithm

The idea behind the first parallel algorithm is to compute the prefix vector in three phases:

- Partitionate the input vector in blocks and compute in parallel the prefix vector of each of them.
- In the second phase we put the final element of each block in a temporary vector and compute its prefix vector.
- Finally for each block i (except for the first one) add in parallel the $(i-1)$ th element of the temporary vector.

At the end of the three phases the final vector contains the correct prefixes.

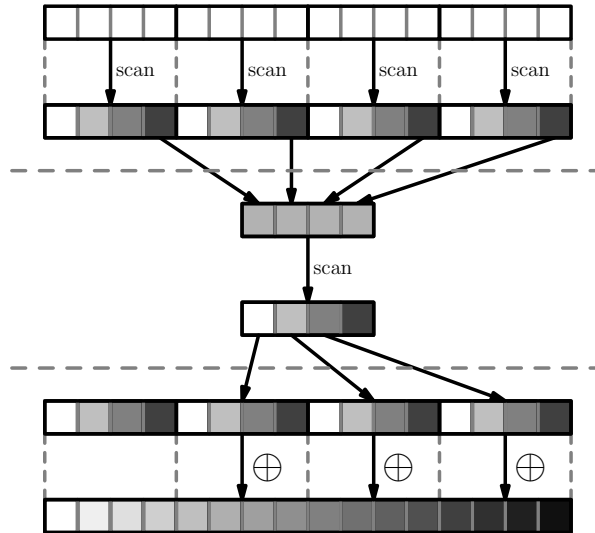


Figure 1: Block-based algorithm graphic representation

3.2 Circuit-based algorithm

Another class of algorithms are the one based on a circuit-like representation. A generic prefix circuit C is a collection of tasks T_1, T_2, \dots, T_t . A task T_i is a set of operations (l, r) that means to compute $x_r = x_r \oplus x_l$. The tasks must be computed sequentially from 1 to t but the advantage is that all the operations in a same task can be performed in parallel without race conditions.

For example the sequential algorithm seen before can be represented as a prefix circuit $S(n) = \{T_1, T_2, \dots, T_{n-1}\}$ where $T_i = \{(i-1, i)\}$.

We present now the simple (but still efficient) prefix circuit $P(2^m)$:

$$P(2^m) = \{T_1, T_2, \dots, T_m, T_{m+1}, \dots, T_{2m-1}\}$$

where for $i = 1, \dots, m$:

$$T_i = \{(2k^i - 2^{i-1} - 1, 2k^i - 1) \mid k = 1, \dots, 2^{m-i}\}$$

and for $i = 1, \dots, m-1$:

$$T_{m+i} = \{(k2^{m-i} - 1, k2^{m-i} + 2^{m-i-1} - 1) \mid k = 1, \dots, 2^i - 1\}$$

In the first phase (from 1 to m) we compute operations like in a reduction tree, in the second phase (from $m+1$ to $2m-1$) we compute the final value of each prefix as the sum of only two others values.

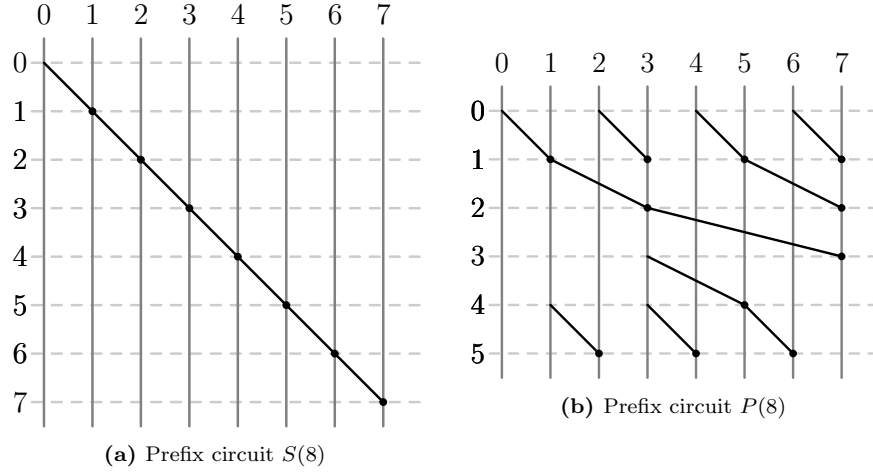


Figure 2: Examples of prefix circuits

As we can see in these simple examples, for $n = 8$ the sequential algorithm needs 7 units of time against the only 5 of the parallel algorithm (obviously if the number of workers is at least 4).

4 Performance modeling

For the theoretical performance modeling we will denote with:

- $n = 2^m$ the size of the input vector ($m = \log_2 n$).
- p the parallelism degree (for simplicity $p \leq n$).
- t_\oplus is the time to compute a single \oplus operation.

4.1 Sequential algorithm

The completion time of the sequential algorithm is simply $T_C^s = (n - 1) \times t_\oplus$.

4.2 Block-based algorithm

The parallel completion time of the block-based algorithm is the sum of the completion times of all of the three phases:

$$T_{C_p}^b \geq T_{1_p} + T_{2_p} + T_{3_p}$$

Where:

- $T_{1_p} = (n/p) \times t_\oplus$ is the time to compute the prefix sums of a single block (as all the p blocks are performed in parallel).
- $T_{2_p} = p \times t_\oplus$ is the time to compute the prefix sums over the temporary vector.
- $T_{3_p} = (n/p) \times t_\oplus$ (or 0 if $p = 1$) is the time to add a value to a single block (except the first one).

So the total completion time is $T_{C_p}^b \geq (2n/p + p) \times t_\oplus$.

In this analysis we assume that n is a multiple of p and all the blocks have the same size, if this is not the case we can consider the block size as the size of the biggest block.

4.3 Circuit-based algorithm

In the prefix circuit $P(n)$ we have a total of $2m - 1$ tasks that must be executed sequentially. As said before all the operations in the same task can be performed in parallel at the same time, so a generic task T_i has a completion time of $1 \times t_\oplus$.

So the completion time of the algorithm is $T_{C_p}^c \geq (2 \log_2(n) - 1) \times t_\oplus$.

However this analysis works only if the parallelism degree is high enough to allow to performs all the operations in T_i at the same time ($p \geq n/2$).

When this is not the case a task T_i is performed in $\max(1, |T_i|/p) \times t_\oplus$ time and the completion time is circa $T_{C_p}^c \simeq 2(\log_2 p + n/p) \times t_\oplus$.

4.4 Comparison table

Now we can compare the completion time of all the algorithms:

p	T_C^s	$T_{C_p}^b$	$T_{C_p}^c$
1	$(n-1) \times t_{\oplus}$	$n \times t_{\oplus}$	$2n \times t_{\oplus}$
2	$(n-1) \times t_{\oplus}$	$n \times t_{\oplus}$	$(n+2) \times t_{\oplus}$
4	$(n-1) \times t_{\oplus}$	$(n/2+4) \times t_{\oplus}$	$(n/2+4) \times t_{\oplus}$
8	$(n-1) \times t_{\oplus}$	$(n/4+8) \times t_{\oplus}$	$(n/4+3) \times t_{\oplus}$
...
$n/2$	$(n-1) \times t_{\oplus}$	$(n/2+1) \times t_{\oplus}$	$2(\log_2(n)+1) \times t_{\oplus}$
n	$(n-1) \times t_{\oplus}$	$(n+2) \times t_{\oplus}$	$2(\log_2(n)+1) \times t_{\oplus}$

Table 1: Comparison of the three algorithms with different parallelism degrees

What we can see from the table is that until $p = 2$ the improvement is negligible, but when the parallelism degree increase the parallel algorithms scales very well.

The circuit-based algorithm reach its best performance for $p \geq n/2$ where each task is computed in $1 \times t_{\oplus}$. The block-based algorithm scales as well as the circuit-based, but when the parallelism degree increase we have a lack of performance due the trade-off between the parallel part of the first and the third stage ($2n/p \times t_{\oplus}$) and the second stage that must be performed sequentially ($p \times t_{\oplus}$).

5 Implementations structure and details

All the implementations are written in C++17, the source code is available as attachment with the report or on GitHub.

Both of the parallel algorithms have been implemented using the standard C++ threads, the FastFlow framework and the OpenMP/Cilk Plus API. All the algorithms are designed to avoid race conditions between the operations in the same phase, so there is no need for any special synchronization mechanism except for the barriers at the end of each phase to pass to the next one.

Each different implementation is available as single class in a single *.hpp* file. All the classes provide the methods to change the parallelism degree, to start the computation and to get time statistics of last computation.

Among with the implementations the project include a benchmark suite that provide different features for the tests: select which algorithm execute, select the range of input size, select the range for the parallelism degree and many other options.

To compile the project is enough to run the *compile.sh* script inside the project folder. If needed a short help is available for consultation, just run *./bin/benchmark -h*.

5.1 Sequential algorithm

The sequential implementation is simply an iteration over the input vector likely the implementation of the *std::partial_sum* function.

5.2 Block-based algorithm

The main problem in the implementation of the block-based algorithm is how to manage, create and assign the blocks to the threads.

To solve this in the project

5.3 Circuit-based algorithm

In the implementation of this algorithm we use again the block manager together with

6 Experimental validation

6.1 experiments details

The experiments were executed on a CPU *Xeon Phi KNC* with 64 physical cores, each of them with a four-way multithreading support (for a total of 256 threads).

- Each running time is expressed in milliseconds and it is an average of 5 repeated test.
- Each implemented algorithm was tested with $p \in \{1, 2, 4, 8, 16, 32, 64, 128\}$.
- Each implemented algorithm was tested with $n \in \{2^{20}, 2^{21}, \dots, 2^{28}, 2^{29}\}$.
- The correctness of each algorithm has been tested with the sequential algorithm.

To replicate the experiments just run the benchmark with these parameters:

```
./benchmark -check -all -m-min 20 -m-max 29 -exp 5 -par-max 128
```

6.2 benchmark results

7 Conclusion

TODO